

PrEstoCloud: Proactive Cloud Resources Management at the Edge for Efficient Real-Time Big Data Processing

Yiannis Verginadis¹, Iyad Alshabani², Gregoris Mentzas¹ and Nenad Stojanovic³

¹ Institute of Communications and Computer Systems, National Technical University of Athens, Athens, Greece

² Activeeon, Valbonne, France, ³ Nissatech Innovation Centre, Nis, Serbia

{jverg, gmentzas}@mail.ntua.gr, iyad.alshabani@activeeon.com, nenad.stojanovic@nissatech.com

Keywords: Big Data Processing, Dynamic Resources Management, Edge Computing

Abstract: Among the greatest challenges of cloud computing is to automatically and efficiently exploit infrastructural resources in a way that minimises cloud fees without compromising the performance of resource demanding cloud applications. In this aspect the consideration of using processing nodes at the edge of the network, increases considerably the complexity of these challenges. PrEstoCloud idea encapsulates a dynamic, distributed, self-adaptive and proactively configurable architecture for processing Big Data streams. In particular, PrEstoCloud aims to combine real-time Big Data, mobile processing and cloud computing research in a unique way that entails proactiveness of cloud resources use and extension of the fog computing paradigm to the extreme edge of the network. The envisioned PrEstoCloud solution is driven by the microservices paradigm and has been structured across five different conceptual layers: i) Meta-management; ii) Control; iii) Cloud infrastructure; iv) Cloud/Edge communication and v) Devices, layers.

1 INTRODUCTION

Cloud computing is “the infrastructure” for real-time data-intensive applications due mainly to its theoretically boundless scalability power. However, in dynamic IoT-based and Big Data driven environments, there are still valid challenges with respect to efficient and reliable real-time processing, since the existing widely adopted technologies (e.g. Hadoop) are not applicable for real-time processing due to their static nature (Mone, 2013). In addition, although more recent and dynamic technologies like STORM (<http://storm.apache.org/>) or SPARK (<http://spark.apache.org/>), offer a very efficient computational solution for soft real-time processing, they don’t exploit new decentralised paradigms (e.g. distributed clouds, edge (Cuomo et al., 2013) and fog computing (Cisco, 2015)).

Regardless of the processing approach that could be adopted, one of the best choices for dealing efficiently with the dynamicity of such domains includes the cloud computing paradigm or extensions of it. In terms of infrastructural support, it is evident that cloud resources scalability is critical for dealing with the most challenging aspects of data-intensive processing. Cloud computing represents a significant shift in the way IT resources have been traditionally

managed and consumed, with significant advantages in terms of cost, flexibility and business agility (Cisco, 2015). Nowadays, the use of cloud computing becomes a necessity in domains characterized by dense real-time sensing and large numbers of distributed heterogeneous information sources. So, it is true that the cloud computing has been recognized as a paradigm for real-time big data processing (Gartner, 2016).

Nevertheless, there are arguments that highlight as a major potential downside of cloud computing the slower-than-desirable performance due to networks’ bandwidth limitations (Mims, 2014). The problem of how to process efficiently on the cloud is becoming all the more acute as more and more objects become “smart,” or able to sense their environments, connect to the Internet, and even receive commands remotely. Modern 3G and 4G cellular networks simply aren’t fast enough to transmit data from devices to the cloud at the pace it is generated. Because of this fact, recently there is a change in focus and new efforts on fog computing become more popular since they store and process the torrent of data being generated by the Internet of Things (IoT) on the “things” themselves, or on devices that sit between our things and the Internet. According to Cisco (2015), the fog extends the cloud to be closer to the things that produce and act on IoT data. These devices, called fog nodes, can

be deployed anywhere with a network connection: on a factory floor, on top of a power pole, alongside a railway track, in a vehicle, or on an oil rig. Any device with computing, storage, and network connectivity can be a fog node. Examples include industrial controllers, switches, routers, embedded servers, and video surveillance cameras.

Nevertheless, the challenge still remains with respect to even higher efficiency and means for dealing with the immense amount of data that the billions of distributed IoT sensors can generate (Shi et al., 2016). Thus, there is the need to extend the fog computing paradigm and reach the extreme edge of the network and beyond, where even data stream sources themselves can participate in aspects of the processing (e.g. smartphones), orchestrated in a controlled way with the rest of cloud resources.

Moreover, the dynamic nature of Big Data environments, which involve real-time changes in data streams, also requires soft real-time (self-) adaptation of cloud resources in order to cope with the imposed scalability challenges. Real-time sources introduce a continuous need for adaptation to the changes in the sensed or observed data, leading to the concept of self-adaptive processing architectures. We note that we focus on soft real-time adaptation of cloud-based systems where the goal is to meet a certain subset of deadlines in order to optimize some application-specific criteria. Thus, a major concern about Information and Communications Technology (ICT) systems operating in a real-time, big data-driven environments, is the anticipation of the reactivity and even the proactiveness to changes in such environment. The inevitable use of cloud resources introduces additional possible failure points and security issues, hence complicating the system adaptation processes that are frequently needed in such dynamic environments and significantly raising costs and security concerns. Moreover, this adaptivity requires managing the usage of cloud resources on another abstraction level, leading to the proactivity in cloud and (more generically) network resource management. Airplanes are a great example of the amount of available data and the benefits of the appropriate processing. In a new Boeing Co. 747, almost every part of the plane is connected to the Internet, recording and, in some cases, sending continuous streams of data about its status. General Electric Co. has said (Mims, 2014) that in a single flight, one of its jet engines generates and transmits half a terabyte of data. Predictive analytics lets such companies know which part of a jet engine might need maintenance, even before the plane carrying it has landed.

Based on this, recent trends in cloud computing go towards the development of new paradigms in the cloud (e.g. heterogeneous, federated, distributed multi-clouds and extending them beyond the fog computing) alleviating the tight interactions between

the computing and networking infrastructures, with the purpose of optimising the use of cloud resources with respect to cost, flexibility and scalability. However, the ever increasing requirements for efficient and resilient data-intensive applications that are able to cope with the variety, volume and velocity of Big Data, lead to the big challenge of new agile architecture paradigms that enhance the dynamic processing even at the extreme edge of the network.

In this paper we describe such a novel processing architecture that deals with such soft real-time adaptation issues, structuring the content as follows: In Section 2, we discuss the envisioned evolution of Real-time Big Data Processing through the PrEstoCloud Approach, providing a detailed conceptual architecture of the solution. In Section 3, we shortly provide the limitations of the current state-of-the-art while in Section 4, we summarise this position paper by discussing the next steps for implementing and validating this work.

2 PRESTOCLOUD CONCEPT

2.1 Evolution of Real-time Big Data Processing through the PrEstoCloud Approach

PrEstoCloud envisions to advance the state-of-the-art with respect to Cloud, Edge computing and real-time data intensive processing in order to provide a dynamic, distributed and proactively configurable architecture for self-adaptive processing of real-time streams. This PrEstoCloud vision can be wrapped around two main drivers:

- in a highly challenging Big Data-driven environment, end users seek for personalised innovative services and superior user-experience that can only be achieved through novel technologies that combine edge analytics, stream mining, processing and exploitation for QoS;
- IT solution providers (esp. SMEs) are facing the limitation of the traditional real-time big data processing architectures, while they need to exploit any business opportunity inherited from dynamic and efficient processing of real-time data streams (incl. on mobile devices) capabilities.

The resolution of these two business needs opens very challenging research questions of designing and developing novel processing architectures based on cloud computing infrastructures and resources at the edge of the network. The main problem in developing data-intensive processing architectures is the dynamicity of the input data which cannot be precisely predicted in advance. We note that we focus

on application scenarios where the speed of an adaptation or a reconfiguration action is not that important as it is to detect or predict a potential issue on time for avoiding system failures and adequate QoS. Contrary to querying historical data, searching in real time streams is related to “guessing” which data will be available in real-time especially with respect to the velocity and variety of data. By assuming that real-time streams are usually coming from sensor devices which can generate a huge amount of data (e.g. 15K events per sec/per sensor) it is clear that dealing with real time streams requires elastic processing infrastructure. Therefore, processing architectures should be able to sense changes in input data streams and adapt themselves accordingly. At the same time the always increasing computing demand will require proactive capability of the resources allocation, i.e. the cloud should be also able to sense situations that require adaptation.

As already discussed, existing real-time processing architectures (like Storm or Lambda (Kiran et al., 2015)), usually don’t exploit multi-cloud environments and present limited adaptation capabilities, while cloud-based data intensive applications, up to now, cannot anticipate adaptations based on real-time changes on the input streams (i.e. proactivity). We summarise the main challenges that we want to point out and address in this position paper as follows: i) exploit multi-cloud environments for deploying Big Data processing frameworks, ii) make intelligent cloud placements and configurations of

applications based on the anticipated processing load with respect to data volume and velocity, iii) elaborate on components that are capable to recommend and implement adaptations in real-time.

PrEstoCloud covers, exactly, these challenges and limitations, contributing to the evolution of real-time Big Data processing. The main driver of the PrEstoCloud architecture is the “change” in data streams, that is an ever emerging concept in real-time data, since it depicts not only velocity (speed of data), but also the “speed of changes” in the Big Data. Traditional monitoring systems are observing when QoS attributes will exceed certain thresholds and then trigger reactions. PrEstoCloud principle is more dynamic, in the sense that it will observe and focus on the dynamics of the trends in changes and reacting if the velocity of the change seems critical (e.g. an increase of the memory consumption on a certain processing node of more than 20% in 30 msec). Therefore, one of the most innovative aspects in PrEstoCloud is the treatment of “changes” as the first class citizen.

2.2 PrEstoCloud Framework

PrEstoCloud aims to address the challenge of cloud-based self-adaptive real-time Big Data processing, including mobile stream processing. The envisioned PrEstoCloud solution is driven by the microservices paradigm that is the use of a software architecture style in which complex applications are composed of

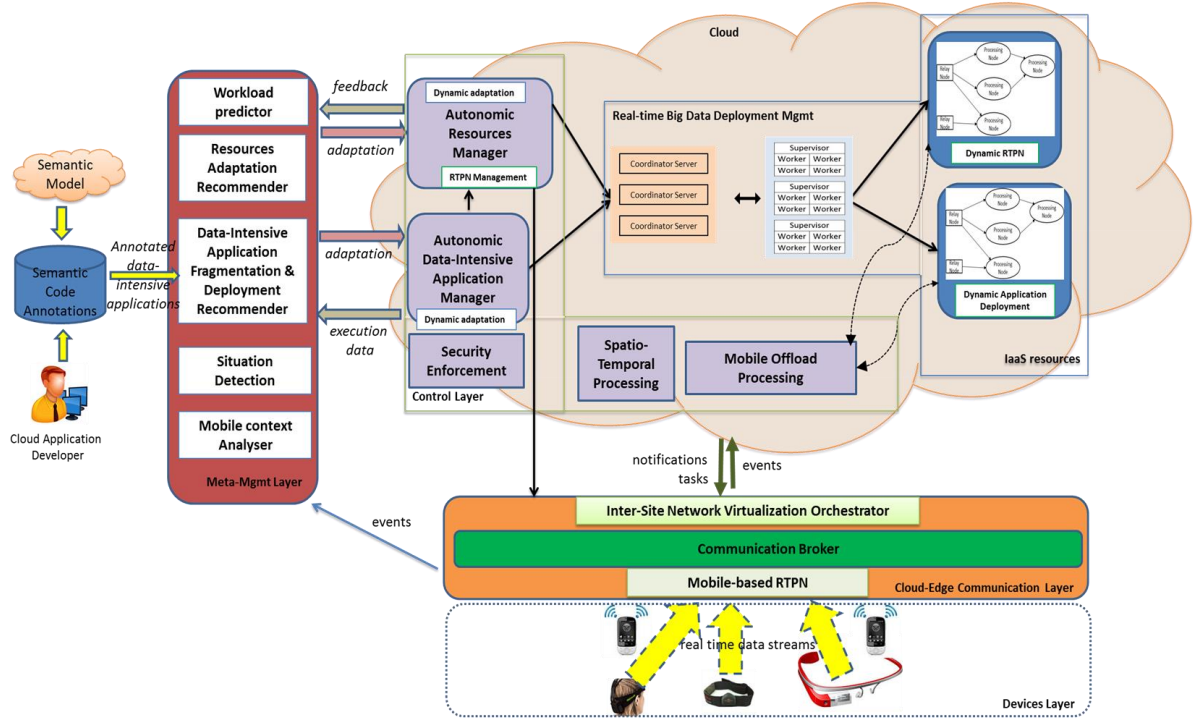


Figure 1: PrEstoCloud Conceptual Architecture

small, independent processes communicating with each other by using language-agnostic APIs. As such, the approach used by PrEstoCloud is superior to any other monolithic solution, since each of the PrEstoCloud components can be instantiated as many times it is needed in order to cope with the scalability requirements of such dynamic environments and satisfy all the availability constraints.

The conceptual architecture of PrEstoCloud is presented in Figure 1 and is briefly discussed below. We note that with respect to real-time data-intensive processing, this solution envisions an architecture inspired by the STORM computational system (i.e. real-time, stateless, stream processing), enhanced in such way that exploits the advantages of multi-cloud environments and Edge computing (i.e. a topology for real-time adaptive and stateful stream processing). In addition, we consider using appropriate models for enabling developers to perform annotations on cloud applications, in order to define their meaningful fragments that may be deployed in a distributed way. The PrEstoCloud architecture has been structured across 5 different layers: i) Meta-management; ii) Control; iii) Cloud infrastructure; iv) Cloud-edge communication and v) Devices layers. The first four layers are discussed below while the fifth one consolidates any kind of device that can be used as a Big Data stream source or as a mobile computational node at the extreme edge of the network.

The *Meta management layer* mainly consists of decision logic capabilities required for enhancing the PrEstoCloud Control layer (e.g. Autonomic Resources Manager, Autonomic Data-Intensive Application Manager). This layer involves the following modules:

Resources Adaptation Recommender that will use as input the situation details, the variation of the Big Data streams and the context of the mobile devices at the extreme edge of the network in order to recommend at the appropriate time, the necessary adaptations of used resources in the real-time processing network (RTPN). These input details will be provided by the Situation Detection mechanism and the Mobile Context Analyser respectively. In addition, these recommendations will be enhanced with the necessary proactiveness based on the envisioned interaction with the Workload Predictor.

Data-Intensive Application Fragmentation & Deployment Recommender that will assist in the appropriate fragmentation of data-intensive applications into smaller parts that can be efficiently deployed over network resources. Specifically, it will be based on the interpretation of the code level annotations (provided by the cloud application developer), in order to recommend meaningful

operations of a cloud application to be distributed and deployed over a STORM-like topology on multi-clouds and Edge resources, for resilience and performance reasons. Different properties like response time, security and locality constraints or other quantitative or qualitative attributes, will be used for deciding on the most optimal configuration at run-time.

Situation Detection Mechanism which will provide the necessary situation awareness for detecting interesting situations by considering Big Data streams and RTPNs deployed on cloud resources or at the extreme edge of the network (i.e. mobile RTPN). Such situations will be completed by any inferred or detected contextual information (provided by the Mobile Context Analyser) and might lead to resources adaptation recommendations or data-intensive application reconfigurations or redeployments. Thus, this mechanism will enhance PrEstoCloud with reliability in terms of the stream processing topology as it will be able to detect possible failures and inform accordingly the two recommenders of the meta-management layer of the PrEstoCloud framework.

Mobile Context Analyser that it will focus on the acquisition and understanding of relevant contextual information derived from any resources at the extreme edge of the network that are or will be engaged in either providing data streams or undertaking parts of the processing effort from the data-intensive applications. The detected or inferred contextual information will be propagated to the Situation Detection Mechanism.

Workload Predictor for fusing proactiveness to the PrEstoCloud solution. Specifically, this module will be able, given an appropriate model and method, recent monitoring information and workload evolution over time to predict the workload that may be experienced in the near future. In addition, it will be also able to predict possible failures because of the overuse of certain processing nodes, thus enhancing the reliability of the processing topology and minimizing any chances for down time incidents.

The second layer of PrEstoCloud architecture is the *Control layer* that manages resources of the *Cloud infrastructure layer* and contains the following modules:

Autonomic Resources Manager which involves monitoring and management of cloud resources capabilities that can be extended to the edge of the network. Specifically, this module is considered as a master-node that involves “Nimbus-like”

(<http://storm.apache.org/>) functionalities, i.e. enhanced nodes with advanced capabilities for reliably (multi-cloud support) and adaptability, dealing additionally with resources management at the edge of the network. This STORM-like approach uses *Coordinator Servers* (i.e. Zookeeper-like functionality for managing stateful nodes across multi-clouds) and *Supervisors* for coordinating the real-time processing network (RTPN) clusters (deployed on the *Cloud infrastructure layer*). This component will interact with Resources Adaptation Recommender and will exploit the spatio-temporal processing capabilities of the respective module. It is composed of the following microservices: i) Smart monitoring, ii) Multi-IaaS manager and connector, iii) Manager of the resources at the edge, iv) Hybrid cloud deployment manager that implements TOSCA (<https://www.oasis-open.org/committees/>) interfaces, v) Resources adaptation manager.

Autonomic Data-Intensive Application Manager which is responsible for the scheduling of big data applications execution. This module also constitutes a “Nimbus-like” node as a master-node which is responsible for distributing meaningful fragments of data-intensive applications across clusters and launching the appropriate workers for controlling the processing workflow (deployed on the *Cloud infrastructure layer*). In this STORM-inspired approach the notion of a *Bolt* is uplifted into a *Processing Node* in the sense that it doesn’t only address data/event processing using small stateless processing units but it may cover more complex operations of a cloud service that should be distributed across different cloud infrastructures for failover and performance issues. This also entails more complex communication and orchestration schemes than just TCP (Darpa, 1981) connections that are traditionally used in bolts taken from current STORM approaches. Thus, we also propose the use of uplifted *Relay Nodes* inspired by STORM *Spouts*. This component will interact with the Autonomic Resources Manager and the Data-Intensive Application Fragmentation & Deployment Recommender. It includes the following microservices: i) Big Data application scheduler which will orchestrate the process execution on the different resources managed by the Autonomic Resources Manager, ii) Processing distribution manager, iii) Application reconfiguration manager

Spatio-Temporal Processing which copes with clustering groups of devices, network congestion detection and geo-fencing capabilities for providing

location-awareness with respect to cloud resources management and on/offloading processing tasks to the devices layer. This actually entails uplifting of raw data (i.e. detect meaningful event patterns) before submitting them to cloud resources for further processing or because of poor network coverage situations, autonomously communicating and processing between devices and sensors in order to transmit the minimum information possible.

Mobile Offload Processing which focuses on shifting parts of processing on the resources at the extreme edge of the network, using local information (e.g. location, power level etc.) as well as information that may be provided by other mobile devices, sensors or cloud resources. It will be instructed accordingly by the Autonomic Resources Manager to look for appropriate mobile devices and efficiently offload processing nodes from the RTPN.

Security enforcement will be responsible for providing the appropriate access and usage control mechanisms with respect to accessing, reallocating and reconfiguring network and edge resources as well as exploiting their processing outcomes.

The four PrEstoCloud layer is the *Cloud-Edge communication* which refers to the following components:

Inter-site network virtualization orchestrator for coping with the need of virtualizing hardware resources situated in multi-cloud environments, managing their orchestration and provisioning across different and heterogeneous providers. This also includes the control of the inter-sites network virtualization process in a secure way, through a close interaction with the Security Enforcement Mechanism. We note that we are examining the provision of a Virtual Network Orchestrator, which will rely on Software Defined Networking (SDN) technology.

Communication Broker which will undertake the responsibility of relaying data streams on and off the PrEstoCloud platform providing standard publish/subscribe event brokering capabilities.

We note that Coordination Servers, Supervisors, Processing nodes and any other PrEstoCloud components can be placed on multi-cloud environments using either VM-based or Container-based deployment for scalability reasons.

3 BEYOND STATE OF THE ART

Among the greatest challenges of cloud computing is to find balance between under-provisioning

infrastructure resources and overprovisioning, which might unnecessarily increase the cloud fees. An alternative way to save resources is to reconfigure the resources allocated to the application, as explored in (Kokkinos et al., 2015) for Amazon EC2. This also requires being able to accurately model the workload variations, which is a tedious task, especially for public clouds where only few trace data are available. A Few recent works paved the way for modelling the dynamicity of the workload experienced in heterogeneous cloud computing platforms (Wolski et al., 2014). In the multi-cloud management side, the main challenges are related to integration issues. Configuration management (CM) tools represent another type deployment management tools focused on the configuration capabilities, which in general address a similar problem. The most representative ones are Chef (<http://chef.io>) and Puppet (<http://puppetlabs.com>). However, they have certain limitations with respect to Fog computing such as the pull-based approach with clients running on the machines or the dependency resolution which is not suitable for resource constrained devices.

Another challenge in the proactive cloud computing is the timely placement of resources on to the infrastructure. To provide efficient placements and evolutive algorithms, complete (but time-bounded) approaches become viable candidates. More specifically, Constraint Programming (CP) (Rossi et al, 2006) based solutions tend to outperform other approaches such as mathematical or logic programming, even more with composite aspects. Multi-Cloud/Fog management and configuration automation extension to the Edge of the network is still an in-progress domain with several unresolved challenges (Shi et al., 2016). Indeed, the Edge part is critical in the sense that it is even more heterogeneous than the cloud resources. PrEstoCloud intends to provide the appropriate extensions to available virtualisation, cluster management and distributed scheduling technologies in order to introduce proactive cloud resources management at the Edge.

4 CONCLUSIONS

This position paper has sketched the PrEstoCloud novel idea which encapsulates a dynamic, distributed, self-adaptive and proactively configurable architecture for processing Big Data streams. Driven by the challenges elaborated in section 2.1 and the limitations of the current approaches, discussed in section 3, this paper suggests the combination of real-time Big Data, mobile processing and cloud computing research in a unique way that entails proactiveness of cloud resources use and extension of the fog computing paradigm to the extreme edge of the network.

The envisioned PrEstoCloud solution is currently under development and its innovative aspects will be thoroughly evaluated through a number of data-intensive pilots. These pilots will stress-test and validate PrEstoCloud claims against real use cases from the transport, journalism and video surveillance domains.

Acknowledgements

This research has received funding from the EU's Horizon 2020 research and innovation programme, under grant agreement No 732339 (PrEstoCloud).

REFERENCES

- Mone, G., 2013. Beyond Hadoop. In *Communications of the ACM*, Vol. 56 (1), pp. 22-24.
- Cuomo, G., Martin, B.,K., Smith, K.,B., Ims, S., Rehn, H., Haberkorn, M., Parikh, J. (2013). Developing Edge Computing Applications. IBM White Paper, Available online at: <http://www.ibm.com/developerworks/cn/websphere/download/pdf/OnDemandEdgeComputing.pdf>
- Cisco (2015) White Paper: Fog Computing and the Internet of Things: Extend the Cloud to Where the Things Are. Available online at: http://www.cisco.com/c/dam/en_us/solutions/trends/iot/docs/computing-overview.pdf
- Darpa, (1981). Transmission Control Protocol. Darpa Internet Program Protocol Specification. Available online at: <https://www.ietf.org/rfc/rfc793.txt>
- Gartner (2016), Innovation Insight for Dynamic Optimization Technology for Infrastructure Resources and Cloud Services. Available online at: https://www.gartner.com/doc/3231420?srcId=1-2819006590&cm_sp=gi_-rr_-top
- Kiran, M., Murphy, P., Monga, I., Dugan, J., Baveja, S., S., 2015. Lambda architecture for cost-effective batch and speed big data processing. In *IEEE International Conference on Big Data (Big Data)*, DOI: 10.1109/BigData.2015.7364082
- Kokkinos, P., Varvarigou, T.A., Kretsis, A., Soumplis, P., Varvarigos E.A., 2015. SuMo: Analysis and Optimization of Amazon EC2 Instances. *J Grid Computing* 13(2): 255-274. doi:10.1007/s10723-014-9311-x
- Mims, M., (2014). Forget 'the Cloud'; 'the Fog' Is Tech's Future. In *The Wall Street Journal*. Available online at: <http://www.wsj.com/articles/SB1000142405270230490830457956662320279406>
- Rossi F., P. van Beek, and T. Walsh, editors. Handbook of Constraint Programming, volume 2 of Foundations of Artificial Intelligence. Elsevier, 2006.
- Shi, W., Cao, J., Zhang, Q., Li, Y., Xu, L., 2016. Edge Computing: Vision and Challenges. *IEEE Internet of Things Journal*, 3(5), pp 637 – 646.
- Wolski, R., Brevik, J., 2014. Using Parametric Models to Represent Private Cloud Workloads. *IEEE Trans. Services Computing* 7(4): 714-725